(54) Title: FACE SYNTHESIS SYSTEM AND METHODOLOGY

(57) Abstract

A system and method for synthesizing a facial image, compares a speech frame from an incoming speech signal with acoustic features stored within visually similar entries in an audio–visual codebook to produce a set of weights. The audio–visual codebook also stores visual features corresponding to the acoustic features. A composite visual feature is generated as a weighted sum of the corresponding visual features, from which the facial image is synthesized. The audio–visual codebook may include multiple samples of the acoustic and visual features for each entry, which corresponds to a sequence of one or more phonemes.

FACE SYNTHESIS SYSTEM AND METHODOLOGY

RELATED APPLICATIONS

This application claims the benefit of U.S. Provisional Application No.

5    60/077,565 entitled "Face Synthesis System and Methodology" filed on March 11,

1998 by Levent M. Arslan and David Talkin, the contents of which are incorporated

by reference herein.

This application contains subject matter related to the commonly assigned, co-

pending PCT patent application No. PCT/US98/01538 entitled "Voice Conversion

10    System and Methodology" filed on January 27, 1998 naming Levent M. Arslan and

David Talkin as inventors, the contents of which are incorporated by reference herein.

FIELD OF THE INVENTION

The present invention relates to audiovisual systems and, more particularly, to

a system and methodology for face synthesis.

15    BACKGROUND OF THE INVENTION

Recently there has been significant interest in face synthesis. Face synthesis

refers to the generation of a facial image in accordance with a speech signal, so that it

appears to a viewer that the facial image is speaking the words uttered in the speech

signal. There are many applications of face synthesis including film dubbing, cartoon

20    character animation, interactive agents, and multimedia entertainment.

Face synthesis generally involves a database of facial images in

correspondence with distinct sounds of a language. Each distinct sound of the

language is referred to as a "phoneme," and during pronunciation of a phoneme, the

mouth and lips of a face form a characteristic, visible configuration, referred to as a

25    "viseme." Typically, the facial image database includes a "codebook" that maps each

phoneme of a language to a corresponding viseme. Accordingly, the input speech text

1

is segmented into phonemes, and the corresponding viseme for each phoneme is sequentially fetched from the database and displayed.

Realistic image quality is an important concern in face synthesis, and transitions from one sound to the next are particularly difficult to implement in a life-

5  like manner because the mouth and lips are moving during the course of pronouncing a sound. In one approach, the mathematical routines are employed to interpolate a series of intermediate images from one viseme at one phoneme to the next. Such an approach, however, can result in an unnatural or distorted appearance, because the movements from one mouth and lip configuration to another are often non-linear.

10  In general, it is practical to store only a restricted number of phoneme/viseme sequences in the codebook. For example, image quality may be improved by storing visemes for all the allophones of a phoneme. An allophone of a phoneme is a slight, non-contrastive variation in pronunciation of the phoneme. A similar issue occurs in applying a face synthesis system originally developed for one language to speech in

15  another language, because the other language includes additional phonemes lacking in the original language. Furthermore, the precise shape of a viseme is often dependent on the neighboring visemes, and there has been some interest in using sequences of phonemes of a given length, such as diphones.

Augmenting the codebook for every possible allophone, foreign phoneme, and

20  phoneme sequences with their corresponding visemes consumes an unacceptably large amount of storage. In a common approach, aliasing techniques are employed in which visemes for a missing phoneme or sequence of phoneme are replaced by existing visemes in the codebook. Aliasing, however, tends to introduce artifacts at the frame boundaries, thereby reducing the realism of the final image.


25  SUMMARY OF THE INVENTION

Accordingly, there exists a need for a face synthesis system and methodology that generates realistic facial images. In particular, there is a need for handling

2

transitions from one viseme to the next with improved realism. Furthermore, a need exists for generating realistic facial images for sequences of phonemes that are missing the codebook or for foreign language phonemes.

These and other needs are addressed by a method and computer-readable
5   medium bearing instructions for synthesizing a facial image, in which a speech frame from an incoming speech signal is compared against acoustic features stored within an audio-visual codebook to produce a set of weights. These weights are used to generate a composite visual feature based on visual features corresponding to the acoustic features, and the composite visual feature is then used to synthesize a facial
10  image. Generating a facial image based on a weighted composition of other images is a flexible approach that allows for more realistic facial images.

For example, more realistic viseme transitions during the course of pronunciation may be realized by using multiple samples of the acoustic and visual features for each entry in the audio-visual codebook, taken during the course of
15  pronouncing a sound. Visemes for foreign phonemes can be generated by combining visemes from a combination of audio-visual codebook entries that correspond to native phonemes. For context-sensitive audio-visual codebooks with a restricted number of phoneme sequences, a weighted combination of features from visually similar phoneme sequences allows for a realistic facial image to be produced for a
20  missing phoneme sequence.

In one embodiment, both the aforementioned aspects are combined so that each entry in the audio-visual codebook corresponds to a phoneme sequence and includes multiple samples of acoustic and visual features. In some embodiments, the acoustic features may be implemented by a set of line spectral frequencies and the
25  visual features by the principal components of a Karhunen-Loewe transform of face points.

Additional objects, advantages, and novel features of the present invention will be set forth in part in the description that follows, and in part, will become

3

apparent upon examination or may be learned by practice of the invention. The

objects and advantages of the invention may be realized and obtained by means of the

instrumentalities and combinations particularly pointed out in the appended claims.


BRIEF DESCRIPTION OF THE DRAWINGS

5          The present invention is illustrated by way of example, and not by way of

limitation, in the figures of the accompanying drawings and in which like reference

numerals refer to similar elements and in which:

           FIG. 1 schematically depicts a computer system that can implement the

present invention;

10          FIGS. 2(a) and 2(b) depict the influence of a modification to the first and

second principal components, respectively, of face point data.

           FIG. 3 depicts a viseme similarity matrix 300 corresponding to phoneme in

American English.

           FIG. 4 is a flowchart illustrating a face synthesis process in accordance with

15    one embodiment of the present invention.


DESCRIPTION OF THE PREFERRED EMBODIMENT

           A method and system for face synthesis are described. In the following

description, for the purposes of explanation, numerous specific details are set forth in

order to provide a thorough understanding of the present invention. It will be

20    apparent, however, to one skilled in the art that the present invention may be practiced

without these specific details. In other instances, well-known structures and devices

are shown in block diagram form in order to avoid unnecessarily obscuring the

present invention.

HARDWARE OVERVIEW

Figure 1 is a block diagram that illustrates a computer system 100 upon which an embodiment of the invention may be implemented. Computer system 100 includes a bus 102 or other communication mechanism for communicating information, and a

5     processor (or a plurality of central processing units working in cooperation) 104 coupled with bus 102 for processing information. Computer system 100 also includes a main memory 106, such as a random access memory (RAM) or other dynamic storage device, coupled to bus 102 for storing information and instructions to be executed by processor 104. Main memory 106 also may be used for storing

10    temporary variables or other intermediate information during execution of instructions to be executed by processor 104. Computer system 100 further includes a read only memory (ROM) 108 or other static storage device coupled to bus 102 for storing static information and instructions for processor 104. A storage device 110, such as a magnetic disk or optical disk, is provided and coupled to bus 102 for storing

15    information and instructions.

Computer system 100 may be coupled via bus 102 to a display 111, such as a cathode ray tube (CRT), for displaying information to a computer user. An input device 113, including alphanumeric and other keys, is coupled to bus 102 for communicating information and command selections to processor 104. Another type

20    of user input device is cursor control 115, such as a mouse, a trackball, or cursor direction keys for communicating direction information and command selections to processor 104 and for controlling cursor movement on display 111. This input device typically has two degrees of freedom in two axes, a first axis (e.g., $x$) and a second axis (e.g., $y$), that allows the device to specify positions in a plane. For audio output

25    and input, computer system 100 may be coupled to a speaker 117 and a microphone 119, respectively.

The invention is related to the use of computer system 100 for face synthesis. According to one embodiment of the invention, face synthesis is provided by

5

computer system 100 in response to processor 104 executing one or more sequences

of one or more instructions contained in main memory 106. Such instructions may be

read into main memory 106 from another computer-readable medium, such as storage

device 110. Execution of the sequences of instructions contained in main memory

5   106 causes processor 104 to perform the process steps described herein. One or more

processors in a multi-processing arrangement may also be employed to execute the

sequences of instructions contained in main memory 106. In alternative

embodiments, hard-wired circuitry may be used in place of or in combination with

software instructions to implement the invention. Thus, embodiments of the

10  invention are not limited to any specific combination of hardware circuitry and

software.

The term "computer-readable medium" as used herein refers to any medium

that participates in providing instructions to processor 104 for execution. Such a

medium may take many forms, including but not limited to, non-volatile media,

15  volatile media, and transmission media. Non-volatile media include, for example,

optical or magnetic disks, such as storage device 110. Volatile media include

dynamic memory, such as main memory 106. Transmission media include coaxial

cables, copper wire and fiber optics, including the wires that comprise bus 102.

Transmission media can also take the form of acoustic or light waves, such as those

20  generated during radio frequency (RF) and infrared (IR) data communications.

Common forms of computer-readable media include, for example, a floppy disk, a

flexible disk, hard disk, magnetic tape, any other magnetic medium, a CD-ROM,

DVD, any other optical medium, punch cards, paper tape, any other physical medium

with patterns of holes, a RAM, a PROM, and EPROM, a FLASH-EPROM, any other

25  memory chip or cartridge, a carrier wave as described hereinafter, or any other

medium from which a computer can read.

Various forms of computer readable media may be involved in carrying one or

more sequences of one or more instructions to processor 104 for execution. For

example, the instructions may initially be borne on a magnetic disk of a remote computer. The remote computer can load the instructions into its dynamic memory and send the instructions over a telephone line using a modem. A modem local to computer system 100 can receive the data on the telephone line and use an infrared

5       transmitter to convert the data to an infrared signal. An infrared detector coupled to bus 102 can receive the data carried in the infrared signal and place the data on bus 102. Bus 102 carries the data to main memory 106, from which processor 104 retrieves and executes the instructions. The instructions received by main memory 106 may optionally be stored on storage device 110 either before or after execution by

10      processor 104.

Computer system 100 also includes a communication interface 120 coupled to bus 102. Communication interface 120 provides a two-way data communication coupling to a network link 121 that is connected to a local network 122. Examples of communication interface 120 include an integrated services digital network (ISDN)

15      card, a modem to provide a data communication connection to a corresponding type of telephone line, and a local area network (LAN) card to provide a data communication connection to a compatible LAN. Wireless links may also be implemented. In any such implementation, communication interface 120 sends and receives electrical, electromagnetic or optical signals that carry digital data streams

20      representing various types of information.

Network link 121 typically provides data communication through one or more networks to other data devices. For example, network link 121 may provide a connection through local network 122 to a host computer 124 or to data equipment operated by an Internet Service Provider (ISP) 126. ISP 126 in turn provides data

25      communication services through the world wide packet data communication network, now commonly referred to as the "Internet" 128. Local network 122 and Internet 128 both use electrical, electromagnetic or optical signals that carry digital data streams. The signals through the various networks and the signals on network link 121 and

7

through communication interface 120, which carry the digital data to and from computer system 100, are exemplary forms of carrier waves transporting the information.

Computer system 100 can send messages and receive data, including program
5   code, through the network(s), network link 121, and communication interface 120. In the Internet example, a server 130 might transmit a requested code for an application program through Internet 128, ISP 126, local network 122 and communication interface 118. In accordance with the invention, one such downloaded application provides for face synthesis as described herein. The received code may be executed
10  by processor 104 as it is received, and/or stored in storage device 110, or other non-volatile storage for later execution. In this manner, computer system 100 may obtain application code in the form of a carrier wave.


### AUDIO-VISUAL CODEBOOK

In accordance with an embodiment of the present invention, an off-line
15  training phase is undertaken as a preliminary step to generate an audio-visual codebook and, preferably, a viseme similarity matrix 300. An audio-visual codebook is data structure that contains entries corresponding to a single phoneme or to a central phoneme in sequence of phonemes, called a "context phone." Each entry includes one or more acoustic features for the phoneme and the corresponding visual features
20  of a related viseme.

The off-line training phase involves collecting data from a test subject by recording the synchronized speech and face point trajectories of the subject. According to one training approach, the subject is asked to utter words, phrases, and sentences for which an orthographic transcription is prepared. The recorded acoustic
25  and visual data are then processed and stored in entries in the audio-visual codebook. The number of entries in the audio-visual codebooks will vary from implementation

to implementation and generally depends on a desired trade-off between face synthesis quality and computational performance.

In one embodiment, the acoustic data is sampled at an appropriate frequency such as 16 kHz and automatically segmented using, for example, a forced alignment

5    to a phonetic translation of the orthographic transcription within an HMM framework using Mel-cepstrum coefficients and delta coefficients as described in more detail in C. Wightman & D. Talkin, *The Aligner User's Manual*, Entropic Reseach Laboratory, Inc., Washington, D.C., 1994. Preferably, the sampled voice data is converted into line spectral frequencies, which can be estimated quite reliably and have a fixed range

10   useful for real-time digital signal processing. The line spectral frequency values for the audio-visual codebook can be obtained by first determining the linear predictive coefficients $a_k$ for the sampled signal according to well-known techniques in the art. For example, specialized hardware, software executing on a general purpose computer or microprocessor, or a combination thereof, can ascertain the linear predictive

15   coefficients by such techniques as square-root or Cholesky decomposition, Levinson-Durbin recursion, and lattice analysis introduced by Itakura and Saito.

In one embodiment, the visual data is obtained as 52 "face points" in three-dimensional space corresponding to points on the subject's face. Since each face point represents $x$, $y$, and $z$ coordinates, the total number of face point parameters is

20   156, thereby constituting a 156-dimensional face point vector. An appropriate transformation technique, such as the Karhunen-Loewe transform, is applied to the face point vector to obtain its principal components. Since the points of a face are highly correlated, a significant reduction in dimensionality can be achieved with only minor distortion. A useful property of using principal components to represent visual

25   features is that the principal components designate the directions that correspond to the most correlated movements. Therefore, modifying the weights of the principal components can be used to animate the underlying face with realistic motions.

For example, the eigenvector with the largest eigenvalue was found to correspond with the movement of the lower jaw, which involves the largest set of correlated points in a speaker's face. Thus, modifying just the first principal component results in moving the lower lip and jaw trajectories. The second principal

5 component was found to correspond with the movement of the sides of the mouth. FIGS. 2(a) and 2(b) depict the effect of adjusting only the first and second principal components of a face point vector, respectively, wherein dark curves represent an original face point trajectory and light curves represent the adjusted face point trajectory.

10 In accordance with one aspect, each phoneme segmented from the speech data is tagged with a "context-phone" symbol indicating the context of phoneme. Specifically, the context-phone symbol indicates the phoneme in the center and one or more neighboring phonemes on either side of the center phoneme in the speech data. For example, the phoneme /eh/ in the word "whenever" has a context-phone symbol

15 of /w_eh_n_eh_v_axr_f/ that includes the three closest neighbors on either side. (The rightmost /f/ phoneme belongs to the following word that begins with an 'f' or 'ph'.) Use of context phones, which form a sequence of phonemes including a center phoneme and neighboring phonemes, allows appropriate context-specific visemes to be generated.

20 In accordance with another aspect, each phoneme in the training data is labeled with multiple, uniformly spaced time locations, for example at five locations, within the course of articulation of the phoneme. The acoustic and visual features, *e.g.* line spectral frequencies and Karhunen-Loewe principal components, are stored . in the audio-visual codebook entry for the phoneme or context-phone. Use of

25 multiple acoustic and visual features allows for a smooth and realistic sequence of visemes to be generated during the course of phoneme articulation.

Thus, the audio-visual codebook includes a number of entries corresponding to a phoneme or a center phoneme and including one or more acoustic features and

one or more corresponding visual features. The audio-visual codebook can be used to

generate facial images by comparing an incoming speech frame with the acoustic

features in the entries to estimate weights for each of the compared acoustic features.

The corresponding visual features are combined as a weighted sum to produce a

5      composite visual feature, which is converted into a facial image. Although

performing this process for all the entries in the audio-visual codebook results in a

very high quality output, it is desirable improve the performance of this process.


## VISEME SIMILARITY MATRIX

In one embodiment, the performance can be significantly improved if phonetic

10     information is known a priori about the incoming speech data being synthesized into

facial images. Specifically, several entries in the audio-visual codebook are selected

whose phonemes or context-phones are most visually similar to the phoneme being

pronounced in each incoming speech frame. Thus, the total number of entries that are

compared with the acoustic feature of the incoming speech frame is reduced to only a

15     few of the most visually similar entries. This selection reduces the computational

overhead of the system and improves the overall performance of the face synthesis

process.

Since, in practice, the training data will not include all possible context-phones

of a given length (or all foreign phonemes and allophones), it is desirable to have

20     some method of associating an unseen context-phone with visually similar entries in

the audio-visual codebook. One visually similar measure is based on the Euclidean

distance of the principal components of face data. This similarity measure can be

automatically generated from the training data and stored in a viseme similarity

matrix 300 by estimating an average principal component vector $m_k$ for each phoneme

25     from the various instances of the phoneme in the training data, as follows:

$$m_k = \frac{1}{T}\sum_{t=1}^{T} p_{kt}, \qquad k \in 1..K, \qquad (1)$$

11

where $K$ represents the total number of phoneme in the language, $T$ represent the total

number of the $k$th phonemes in the training data, and $p_{kt}$ represents the $t$th principal

component vector that is associated with the $k$th phoneme. Given the average

principal component vectors $m_k$, the Euclidean distance between each pair of

5      phonemes is calculated as:

$$d_{ik} = \|\mathbf{m}_i - \mathbf{m}_k\|, \qquad i, k \in 1..K, \tag{2}$$

Based on the calculated Euclidean distances, the viseme similarity measure $s_{ik}$

is derived as follows:

$$s_{ik} = e^{-50d_{ik}}, \qquad i, k \in 1..K, \tag{3}$$

10     One property of this formulation is that viseme similarity values $s_{ik}$ will range

between 0 and 1. FIG. 3 depicts a gray scale image of one viseme similarity matrix

300 corresponding to phonemes of American English, wherein darker points represent

a higher level of visual similarity. For example, the most visually similar phoneme to

/b/ was identified to be /p/. In general, it has been found that entries in the viseme

15     similarity matrix 300 agree with intuitive expectations.

The viseme similarity matrix 300 can be used directly to determine the visual

similarity of two phonemes, but a more involved procedure is used to estimate a

visual similarity measure between two context-phones representing a sequence of

phonemes. Preferably, the center phoneme should have the highest influences with

decreasing influence for the phonemes that are more remote from the center phoneme.

One procedure to estimate the visual similarity of context-phones can be formulated

as follows:

$$v_j = s_{cj} \sum_{l=1}^{C} 10^{-i}(s_{lij} + s_{rij}), \quad j \in 1..L, \tag{4}$$

where $C$ is the level context information (*i.e.* the number of neighboring phonemes on

each side), $L$ is the total number of content-phones in the audio-visual codebook, $s_{lij}$ is

the similarity between the $i$th left phoneme of the subject context-phone and the $j$th

20     context-phone in the audio-visual codebook, and $s_{rij}$ is the similarity between the $i$th

right phoneme of the subject context-phone and the $j$th context-phone in the audio-

codebook was configured to store entries for only a single phone, then the viseme similarity matrix 300 is consulted directly to obtain the visual similarity measure.

Based on the determined visual similarity measures, the N most visually similar entries of the audio-visual codebook are selected. The best value for N will vary from implementation to implementation, depending on factors such as the length of the phoneme sequence of the context-phones and the desired performance/realism tradeoff for a given set of training data. Generally, however, the value of N ranges from about four to about sixteen, and may in fact be a user-configurable parameter.

At step 402, the incoming speech frame is converted into acoustic features suitable for comparison with the acoustic features stored in the audio-visual codebook. For example, the incoming speech frame may be converted into line spectral frequencies and compared with a line spectral frequencies set stored in the audio-visual codebook. In some embodiment, a plurality of samples, such as five samples, are stored for each entry in the audio-visual codebook. The result of the acoustic feature comparison is a weight, wherein a higher weight is assigned for more acoustically similar samples. A variety of techniques for producing the weight based on the comparison may be employed but the present invention is not limited to any particular weight.

One weighting technique is described in the commonly assigned, co-pending PCT patent application No. PCT/US98/01538 entitled "Voice Conversion System and Methodology." As described therein, codebook weights $v_i$ are estimated by comparing the input line spectral frequency vector $w_k$ with each acoustic feature sample, $S_i$ in the audio-visual codebook to calculate a corresponding distance $d_i$:

$$d_i = \sum_{k=1}^{P} \mathbf{h}_k \mid \mathbf{w}_k - \mathbf{S}_{ik} \mid, i \in 1..L \tag{5}$$

where $L$ is the codebook size. The distance calculation may include a weight factor $\mathbf{h}_k$, which is based on a perceptual criterion wherein closely spaced line spectral frequency pairs, which are likely to correspond to formant locations, are assigned higher weights:

14

$$\mathbf{h}_k = \frac{e^{-0.05 \cdot |K - k|}}{\min(|\mathbf{w}_k - \mathbf{w}_{k-1}|, |\mathbf{w}_k - \mathbf{w}_{k+1}|)}, k \in 1..P \tag{6}$$

where $K$ is 3 for voiced sounds and 6 for unvoiced, since the average energy decreases
(for voiced sounds) and increases (for unvoiced sounds) with increasing frequency.
Based on the calculated distances $d_i$, the normalized codebook weights $\mathbf{v}_i$ are obtained
as follows:

$$\mathbf{v}_i = \frac{e^{-\gamma d_i}}{\sum_{l=1}^{L} e^{-\gamma d_l}}, i \in 1..L \tag{7}$$

where the value of $\gamma$ for each frame is found by an incremental search in the range of
0.2 to 2.0 with the criterion of minimizing the perceptual weighted distance between
the approximated line spectral frequency vector $\mathbf{v}S_k$ and the input line spectral
frequency vector $\mathbf{w}_k$. These weights may be further adjusted as also described in
pending PCT patent application No. PCT/US98/01538.

At step 404, a composite visual feature is constructed from the weights and the
corresponding visual features of the selected audio-visual codebook entries, for
example, as a weighted sum or linear combination of the principal components of the
facial data samples. For example, the composite visual feature may be calculated as
follows:

$$\tilde{\mathbf{p}} = \sum_{n=1}^{SN} \mathbf{v}_n \mathbf{p}_n. \tag{8}$$

In one embodiment, a plurality of visual features is stored for each entry in the
audio-visual codebook at different points in time during the articulation of the sound
corresponding to the entry. Thus, the weighted sum will include all the visual
samples for the audio-visual codebook entry, thereby producing facial data that more
realistically tracks the movement of the mouth and lips during speaking.

At step 406, the composite visual feature is converted into the desired facial
data. For example, if principal components obtained from a Karhunen-Loewe
transformation are used to represent the visual features, then an inverse Karhunen-
Loewe transformation is applied on the composite principal components to produce

15

face points as output.  These face points can be converted to the facial image by known techniques.

5      Accordingly, a face synthesis system and methodology is described wherein realistic facial images are produced in accordance with an input speech signal. Specifically, a composite visual feature is generated from entries in an audio-visual codebook according to weights identified by comparing the incoming acoustic features with the audio-visual codebook acoustic features.  Consequently, realistic output is attained for viseme transitions, for highly context-dependent situations, and

10     even for foreign language phonemes without requiring the audio-visual codebook to store an enormous amount of training samples.

While this invention has been described in connection with what is presently considered to be the most practical and preferred embodiment, it is to be understood that the invention is not limited to the disclosed embodiment, but on the contrary, is

15     intended to cover various modifications and equivalent arrangements included within the spirit and scope of the appended claims.

CLAIMS

WHAT IS CLAIMED IS:

1        1. A method of synthesizing a facial image in accordance with a speech signal,

2    comprising the steps of:

3        comparing a speech frame of the speech signal with a plurality of acoustic features

4            within an audio-visual codebook to produce therefrom a plurality of weights;

5        generating a composite visual feature based on the weights and a plurality of

6            visual features corresponding to the acoustic features; and

7        synthesizing the facial image based on the composite visual feature.


1        2. The method of claim 1, wherein:

2    the audio-visual codebook contains entries each including a plurality of acoustic

3            features and a plurality of corresponding visual features; and

4    comparing the speech frame includes comparing the speech frame with the

5            acoustic features from an entry within the audio-visual codebook.


6        3. The method of claim 1, wherein:

7    the audio-visual codebook contains entries each including an acoustic feature and

8            a corresponding visual feature; and

9    comparing the speech frame includes comparing the speech frame with acoustic

10           features from a plurality of selected entries within the audio-visual codebook.


11       4. The method of claim 3, wherein the speech signal includes a sequence of

12   speech frames correlated with a sequence of phonemes, said method further

13   comprising:

17

14      determining a visual similarity measure between a phoneme in the sequence that is

15          correlated to the speech frame and the entries in the audio-visual codebook,

16          said entries in the audio-visual codebook corresponding to a phoneme; and

17      selecting the selected entries from the audio-visual codebook based on the visual

18          similarity measure; and

1      5. The method of claim 3, wherein the speech signal includes a sequence of

2   speech frames correlated with a sequence of phonemes, said method further

3   comprising:

4      determining visual similarity measures between a phoneme in the sequence with

5          neighboring phonemes thereof and the entries in the audio-visual codebook,

6          said entries in the audio-visual codebook corresponding to a series of

7          phonemes having a center phoneme with neighboring phonemes thereof; and

8      selecting the selected entries in the audio-visual codebook based on the

9          determined visual similarities.

1      6. The method of claim 5, wherein determining the visual similarity measures

2   includes calculating Euclidean distances between each of sets of principal components

3   of facial data corresponding to the phoneme in the sequence with the neighboring

4   phonemes thereof and principle component samples of facial data corresponding to

5   the center phoneme with the neighboring phonemes thereof.

1      7. The method of claim 5, wherein determining the visual similarity measures

2   includes accessing a visual similarity matrix containing elements based on Euclidean

3   distances between each of sets of principal components of facial data corresponding to

4   the phoneme in the sequence with the neighboring phonemes thereof and principle

5   component samples of facial data corresponding to the center phoneme with the

6   neighboring phonemes thereof.

1    8. The method of claim 1, wherein the acoustic feature includes a line spectral

2    frequencies set and the visual feature includes a set of principal components of facial

3    data derived from face point samples.


1    9. A method of synthesizing a facial image in accordance with a speech signal,

2    said speech signal including a sequence of speech frames correlated with a sequence

3    of phonemes, said method comprising the steps of:

4    determining visual similarity measures between a phoneme in the sequence with

5            neighboring phonemes thereof and entries in an audio-visual codebook, said

6            entries in the audio-visual codebook corresponding to a series of phonemes

7            having a center phoneme with neighboring phonemes thereof and including a

8            plurality of acoustic features and a plurality of corresponding visual features;

9    selecting a plurality of the entries in the audio-visual codebook based on the

10           determined visual similarities;

11   comparing a speech frame of the speech signal with the acoustic features from

12           entries to produce therefrom a plurality of weights;

13   generating a composite visual feature based on the visual features of the entries

14           and the weights; and

15   synthesizing the facial image based on the composite visual feature.


1    10. A computer-readable medium bearing instructions for synthesizing a facial

2    image in accordance with a speech signal, said instructions arranged, when executed

3    by one or more processors, to cause the one or more processors to perform the steps

4    of:

5            comparing a speech frame of the speech signal with a plurality of acoustic features

6                    within an audio-visual codebook to produce therefrom a plurality of weights;

7          generating a composite visual feature based on the weights and a plurality of

8             visual features corresponding to the acoustic features; and

9          synthesizing the facial image based on the composite visual feature.

1          11. A computer-readable medium bearing instructions for synthesizing a facial

2   image in accordance with a speech signal, said speech signal including a sequence of

3   speech frames correlated with a sequence of phonemes, said instructions arranged,

4   when executed by one or more processors, to cause the one or more processors to

5   perform the steps of:

6          determining visual similarity measures between a phoneme in the sequence with

7             neighboring phonemes thereof and entries in an audio-visual codebook, said

8             entries in the audio-visual codebook corresponding to a series of phonemes

9             having a center phoneme with neighboring phonemes thereof and including a

10            plurality of acoustic features and a plurality of corresponding visual features;

11          selecting a plurality of the entries in the audio-visual codebook based on the

12             determined visual similarities;

13          comparing a speech frame of the speech signal with the acoustic features from

14             entries to produce therefrom a plurality of weights;

15          generating a composite visual feature based on the visual features of the entries

16             and the weights; and

17          synthesizing the facial image based on the composite visual feature.
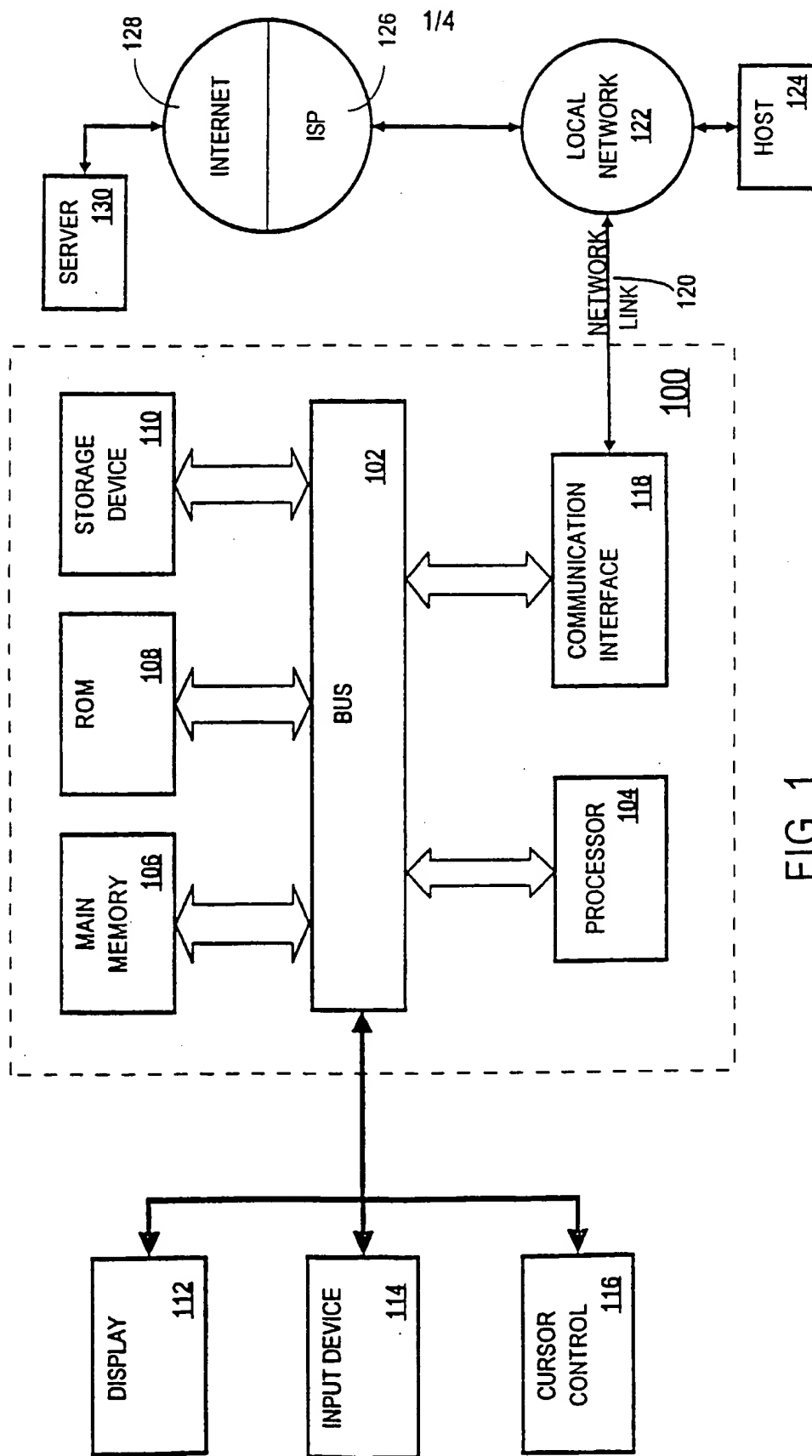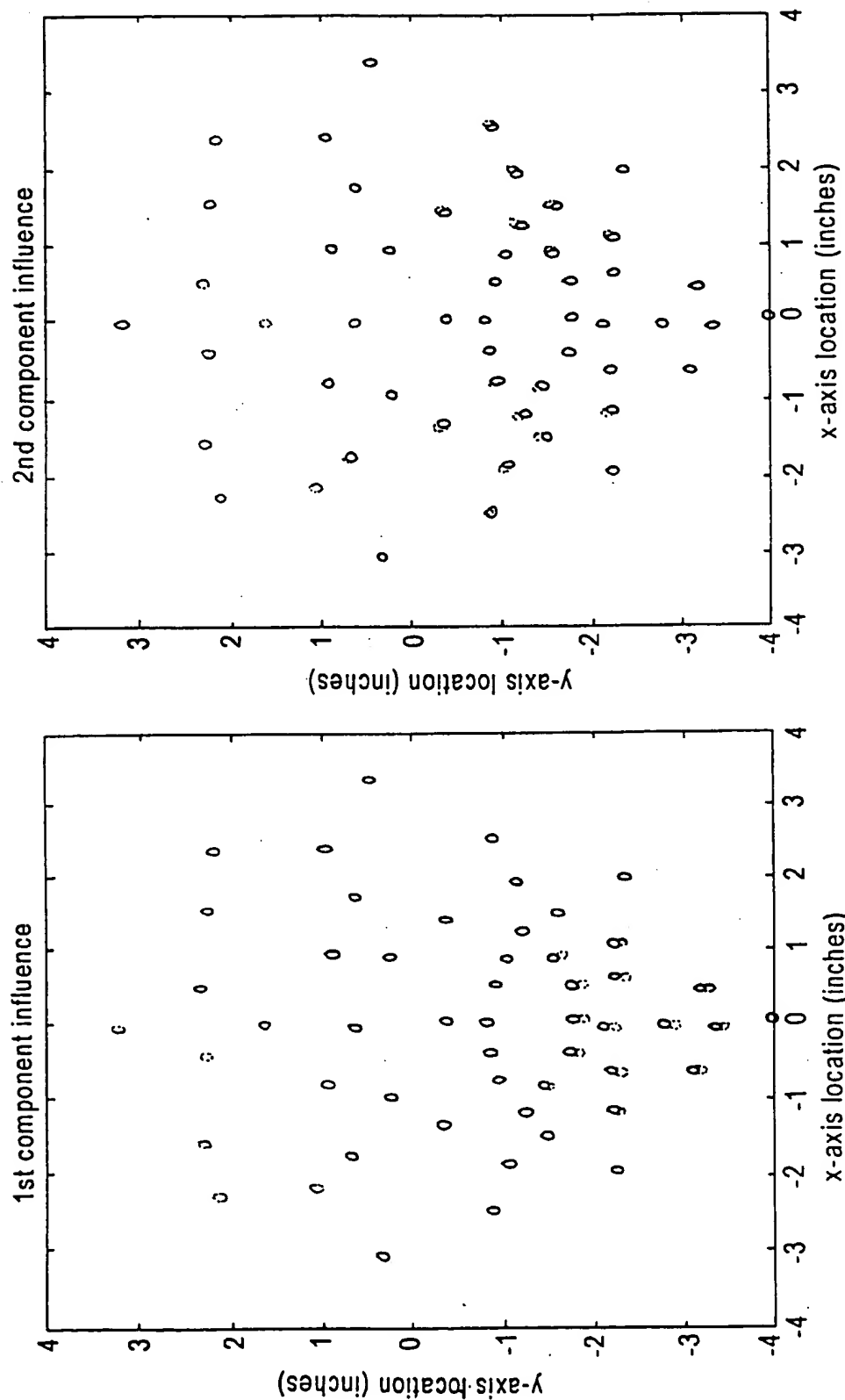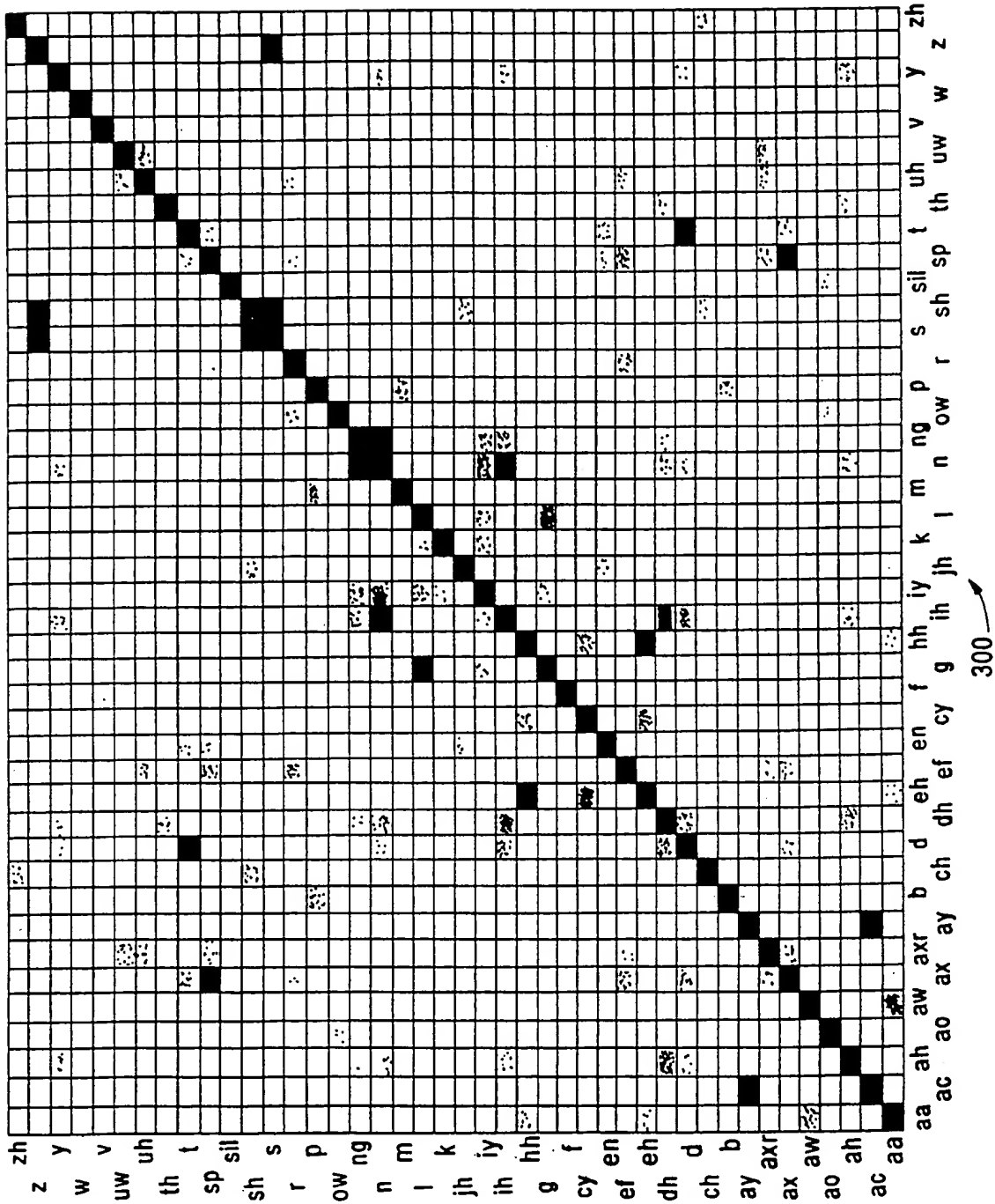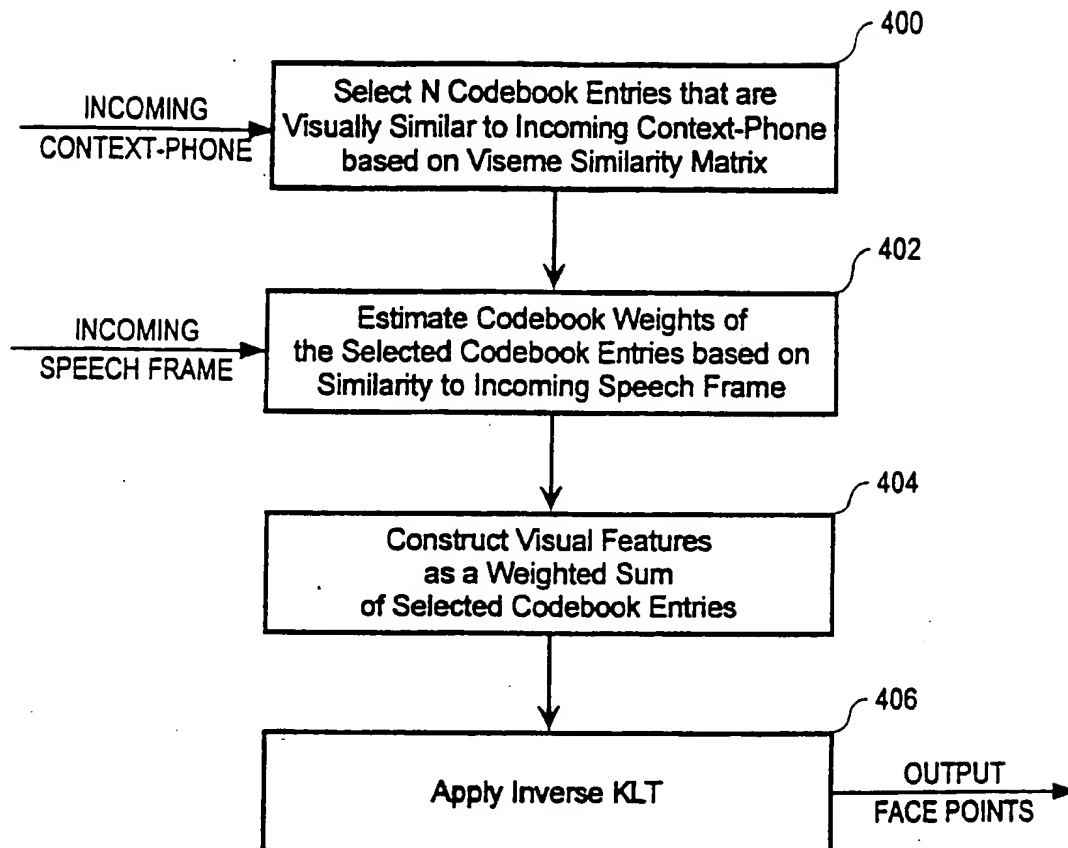
1/4



FIG. 1

2/4



FIG. 2(b)

FIG. 2(a)

3/4



FIG. 3

4/4

```
INCOMING          ┌─────────────────────────────────┐ ⌐ 400
CONTEXT-PHONE ───▶ │  Select N Codebook Entries that are │
                   │ Visually Similar to Incoming Context-Phone │
                   │  based on Viseme Similarity Matrix │
                   └─────────────────────────────────┘
                                    │
                                    ▼
INCOMING          ┌─────────────────────────────────┐ ⌐ 402
SPEECH FRAME ───▶ │   Estimate Codebook Weights of   │
                   │ the Selected Codebook Entries based on │
                   │  Similarity to Incoming Speech Frame │
                   └─────────────────────────────────┘
                                    │
                                    ▼
                   ┌─────────────────────────────────┐ ⌐ 404
                   │     Construct Visual Features    │
                   │        as a Weighted Sum         │
                   │    of Selected Codebook Entries  │
                   └─────────────────────────────────┘
                                    │
                                    ▼
                   ┌─────────────────────────────────┐ ⌐ 406
                   │                                  │   OUTPUT
                   │         Apply Inverse KLT        │ ──▶ FACE POINTS
                   │                                  │
                   └─────────────────────────────────┘
```

# FIG. 4

# INTERNATIONAL SEARCH REPORT

## A. CLASSIFICATION OF SUBJECT MATTER
IPC 6   G06T15/70

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
IPC 6   G06T

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|---|---|---|
| X | EP 0 710 929 A (AT & T CORP) 8 May 1996 (1996-05-08) | 1-5,9-11 |
| Y | abstract; claims 1,7 column 5, line 41 - line 44 | 6-8 |
| Y | US 4 907 276 A (ALDERSBERG SHABTAI) 6 March 1990 (1990-03-06) abstract; claims 1-3 | 6-8 |
| A | US 5 630 017 A (MATTHEWS III JOSEPH H  ET AL) 13 May 1997 (1997-05-13) column 38, line 3 - line 17 claim 10 | 1-11 |
| A | GB 2 231 246 A (KOKUSAI DENSHIN DENWA CO LTD) 7 November 1990 (1990-11-07) page 4, line 19 - page 5, line 21; claims 1-3 | 1-11 |

-/--

| X | Further documents are listed in the continuation of box C. |
| X | Patent family members are listed in annex. |

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

| Date of the actual completion of the international search | Date of mailing of the international search report |
|---|---|
| 21 July 1999 | 30/07/1999 |

| Name and mailing address of the ISA | Authorized officer |
|---|---|
| European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Tx. 31 651 epo nl, Fax: (+31-70) 340-3016 | Filloy García, E |

1

| C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT | | |
|---|---|---|
| Category * | Citation of document, with indication,where appropriate, of the relevant passages | Relevant to claim No. |
| A | WO 97 36288 A (BREEN ANDREW PAUL ;BOWERS EMMA JANE (GB); BRITISH TELECOMM (GB)) 2 October 1997 (1997-10-02) page 11, paragraph 1 | 5-7 |
| A | SIROVICH L ET AL: "LOW-DIMENSIONAL PROCEDURE FOR THE CHARACTERIZATION OF HUMAN FACES" JOURNAL OF THE OPTICAL SOCIETY OF AMERICA - A, vol. 4, no. 3, 1 March 1987 (1987-03-01), pages 519-524, XP000522491 ISSN: 0740-3232 page 1, left-hand column, paragraph 2 | 6-8 |
| A | EP 0 689 362 A (AT & T CORP) 27 December 1995 (1995-12-27) column 4, line 47 - line 51 | 1-11 |
| A | CHOU W ET AL: "SPEECH RECOGNITION FOR IMAGE ANIMATION AND CODING" PROCEEDINGS OF THE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH, AND SIGNAL PROCESSING (ICASSP), DETROIT, MAY 9 - 12, 1995 IMAGE AND MULTI-DIMENSIONAL SIGNAL PROCESSING/ SIGNAL PROCESSING APPLICATIONS DEVELOPMENT, vol. 4, no. CONF. 20, 9 May 1995 (1995-05-09), pages 2253-2256, XP000535403 INSTITUTE OF ELECTRICAL AND ELECTRONICS ENGINEERSISBN: 0-7803-2432-3 | |

1

# INTERNATIONAL SEARCH REPORT

Information on patent family members

| Patent document cited in search report | | Publication date | Patent family member(s) | | Publication date |
|---|---|---|---|---|---|
| EP 0710929 | A | 08-05-1996 | AU | 3668095 A | 16-05-1996 |
| | | | CA | 2162199 A | 08-05-1996 |
| | | | JP | 8235384 A | 13-09-1996 |
| US 4907276 | A | 06-03-1990 | NONE | | |
| US 5630017 | A | 13-05-1997 | US | 5689618 A | 18-11-1997 |
| | | | US | 5613056 A | 18-03-1997 |
| GB 2231246 | A | 07-11-1990 | JP | 2234285 A | 17-09-1990 |
| | | | JP | 2518683 B | 24-07-1996 |
| WO 9736288 | A | 02-10-1997 | AU | 2167097 A | 17-10-1997 |
| | | | CA | 2249016 A | 02-10-1997 |
| | | | CN | 1214784 A | 21-04-1999 |
| | | | EP | 0890168 A | 13-01-1999 |
| EP 0689362 | A | 27-12-1995 | US | 5608839 A | 04-03-1997 |
| | | | CA | 2149068 A | 22-12-1995 |
| | | | JP | 8023530 A | 23-01-1996 |